

# MENGYU YE

Ph.D. Student in Information Science | LLM Interpretability & Multilingual AI

Tohoku University, Sendai, Japan

Email: ye.mengyu.s1@dc.tohoku.ac.jp

URL: <https://muyo8692.github.io>

## Education

---

<b>Tohoku University</b>	Sendai, Japan
<i>Ph.D. Student in Information Science</i>	<i>April 2024 - 2027 (expected)</i>
Advisor: Prof. Jun Suzuki	
<i>M.S. in Information Science</i>	<i>April 2022 - 2024</i>
Advisor: Prof. Jun Suzuki	
<i>B.S. in Engineering</i>	<i>April 2018 - 2022</i>
Advisor: Prof. Xiao Zhou	

## Awards & Grants

---

Gemma 2 Academic Research Program Grant (JP/KR 2024)	2024
BOOST Fellowship	2024
Best Paper Award – ACL 2023 Student Research Workshop	2023

## Publications

- 
- Mengyu Ye**, Tatsuki Kuribayashi, Jun Suzuki, Goro Kobayashi, Hiroaki Funayama. Assessing Step-by-Step Reasoning Against Lexical Negation: A Case Study on Syllogism. *In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP 2023)*. [pdf]
    - Mengyu Ye**, Tatsuki Kuribayashi, Jun Suzuki, Goro Kobayashi, Hiroaki Funayama. Assessing Chain-of-Thought Reasoning Against Lexical Negation: A Case Study on Syllogism. *Non-archival submission for ACL-SRW 2023*.  
🏆 Best Paper Award at ACL-SRW 2023
  - Hiroto Kurita, Ikumi Ito, Hiroaki Funayama, Shota Sasaki, Shoji Moriya, **Ye Mengyu**, et al. TohokuNLP at SemEval-2023 Task 5: Clickbait Spoiling via Simple Seq2Seq Generation and Ensembling. *In Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)*. [pdf]

## Skills

---

*Programming Languages:* Python, MATLAB, Java, C/C++, SQL (MySQL), SML#

*Languages:* Languages: Chinese (Native), Japanese (Near-Native), English (Professional Working Proficiency)

*Machine Learning Frameworks:* JAX/Flax, Pytorch

*Tools & Methods:* ML pipeline development, High-Performance Computing (HPC) Clusters (ABCI, mdx, Genkai etc.), Containerization (Docker, Singularity etc.), Synthetic data generation

## Experience

- 
- Moonshot Research and Development Program** *Sept. 2022 - Present*  
*Research Assistant*
- Conducted fine-tuning of LLaMA models to develop advanced Japanese LLMs, enhancing their linguistic and contextual performance.
  - Designed and implemented a RAG system tailored for integration into prototype cybernetic avatar applications.

## **Tohoku University**

*April 2023 - Present*

### *Research Assistant*

- Designed, implemented machine learning pipelines for large-scale model evaluation, involving dozens of LLMs, multiple tasks, feature analysis while ensuring efficient, reproducible workflows.
- Developed an end-to-end pipeline evaluating input attribution methods in LLMs within ICL settings through synthetic task generation, LLM fine-tuning, task accuracy assessment, attribution performance measurement.
- Collaborated with international teams to create multilingual evaluation frameworks for LLMs, ensuring cross-lingual consistency and robust evaluation metrics.

## **Projects**

---

### **Training Sparse Autoencoders for Japanese LLM**

- Training and plan to releasing Sparse Autoencoders trained on Japanese LLMs to contribute valuable resources to the research community.

### **Cross-Cultural AI Safety Initiative**

- Collaborating with Georgia Tech researchers to develop evaluation frameworks for cultural bias detection in Asian languages.
- Leading data creation efforts for Japanese and Chinese, ensuring cross-lingual consistency.

### **Mechanistic Interpretability for Multilingual Foundation Models**

- Exploring techniques to enhance reliability and explainability in multilingual AI systems, currently in the conceptual stage.

## **Teaching Experience**

---

### **Teaching Assistant**

2024 Spring - Seminar on System Information Sciences

2024 All year - Advanced Seminar on System Information Sciences B

### **Mentorship**

2023 - Kazuki Yano, master's student researcher from Tohoku university GSIS department

2024 - Koichi Iwakawa, master's student researcher from Tohoku university GSIS department

2024 - Haochen Zhu, master's student researcher from Tohoku university GSIS department

## **Service**

---

2024 - ACL Rolling Review (ARR)